

COMANDI STATISTICA

❖ Statistica descrittiva

NOME	COMANDO
Cancella tutto	<code>rm(list=ls())</code>
Vettore	<code>vettore<-c()</code>
Lunghezza	<code>str()</code> → num [1:10] <code>length()</code> → 10
Dati in ordine crescente	<code>sort()</code>
Frequenza assoluta (moda)	<code>table(a)</code>
Frequenza relativa	<code>prop.table(b)</code>
Frequenza percentuale	<code>prop.table(b)*100</code>
Media, xbar	<code>mean()</code>
Mediana	<code>median()</code>
Media ponderata, wa	<code>sum(a*freq_a)/tot_a</code> con <code>tot_a = sum(freq_a)</code>
Deviazione standard, s	<code>sd</code> o <code>sqrt(var)</code>
Varianza, s ²	<code>var()</code>
Coefficiente di variazione, cv	<code>sd()/mean()</code> con mean in valore assoluto
Boxplot	<code>boxplot(a, horizontal=TRUE, col="yellow")</code> <code>boxplot(a,b, horizontal=TRUE,</code> <code>main="titolo",</code> <code>names=c("a","b"),</code> <code>col=c("orange","lightblue"))</code>
Scatterplot → length uguale	<code>plot(a,b,col="red")</code>
Confronto tra percentili → length non uguale	<code>qqplot(a,b,col="red")</code>
Istogramma	<code>hist()</code> breaks=numero, è il n° di classi
Diagramma a torta	<code>pie(tabella frequenze)</code>
Diagramma a barre	<code>barplot(tabella frequenze)</code>
Min, 1st qu., Mediam, Mean, 3rd qu., Max	<code>summary()</code>
Percentile	<code>quantile(vettore,k%)</code>
Covarianza	<code>cov(x,y)</code>
Coefficiente di corr di Pearson, r	<code>cor(x,y)</code> devo avere sd() e cov()
Retta di regressione	<code>lm(y~x)</code> con y=output e x=input <code>y<-function(x) (m*x+q)</code>
Sovrapporre	<code>abline()</code>
Valori attesi	<code>attesi<-round(predict(reg), digits=0)</code> <code>setNames(attesi, vettore richiesto)</code>
Per agire con un comando sulla singola variabile del dataframe	<code>comando(dataframe \$ singola variabile)</code>
Estremi del range	<code>range(dataframe \$ singola variabile)</code>
Più piccolo e più grande valore di tutta la tabella	<code>range(dataframe)</code>
Ampiezza del range	<code>range(dataframe \$ singola variabile)[2]-</code> <code>range(dataframe \$ singola variabile)[1]</code>
Intervallo interquartile (distanza tra q ₃ -q ₁)	<code>IQR(dataframe \$ singola variabile)</code>
Parametri di centralità per tutte le coppie di variabili	<code>summary(dataframe)</code>
Parametri di dispersione per tutte le coppie di variabili (covarianza e coeff. Di Pearson)	<code>cov(dataframe)</code> <code>cor(dataframe)</code>

❖ Probabilità

$P(A)$	k/n
Unione	$A \cup B = \{x \in \Omega; x \in A \text{ o } x \in B\}$
Intersezione	$A \cap B = \{x \in \Omega; x \in A \text{ e } x \in B\}$
Complementare	$A^c = \{x \in \Omega; x \notin A\}$ $(A \cup A^c) = \Omega$ $(A \cap A^c) = 0$
Assiomi di Kolmogorov Spazio di probabilità (Ω, a, P)	<ul style="list-style-type: none"> $0 \leq P(A) \leq 1$ $P(\Omega) = 1$ $A \cap B = 0 \rightarrow P(A \cup B) = P(A) + P(B)$
1. Probabilità del complementare	$P(A^c) = 1 - P(A)$
2. Probabilità dell'evento impossibile	$P(0) = 0$
3. Probabilità dell'evento certo	$P(B) = P(B \cap A) + P(B \cap A^c)$
4. Probabilità di ordinamento	$A \subset B \rightarrow P(A) \leq P(B)$
5. Probabilità dell'unione di eventi non disgiunti	$P(A \cup B) = P(A) + P(B) - P(A \cap B)$
5bis. Estensione punto 5	$P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - (B \cap C) + P(A \cap B \cap C)$
Eventi indipendenti	$P(A \cap B) = P(A) * P(B)$ solo se A e B sono indipendenti
Spazio di probabilità uniforme	$P(A) = P(A) * p = A / \Omega $
	$P(A \cap B) \leq P(A) \leq P(A \cup B)$
Valore atteso=media, μ , $E[x]$	$\sum(x * x)$ $n * p$
Varianza, σ^2	$\sum(x - \mu)^2 * x$ $n * p(1 - p)$
Deviazione standard, σ	$\sqrt{\sigma^2}$
$n!$	$\text{factorial}(n)$
$\binom{n}{k}$	$\text{choose}(n, k)$

$P(x=k)$	$p(k) = P(x=k) = \binom{n}{k} p^k (1-p)^{n-k}$ $\text{dbinom}(k, \text{size}=n, \text{prop}=p)$
$P(a \leq x \leq b)$	$\text{sum}(\text{dbinom}(a:b, \text{size}=n, \text{prop}=p))$
$P(a \leq x < b)$	$\text{sum}(\text{dbinom}(a:b-1, \text{size}=n, \text{prop}=p))$
$P(a < x \leq b)$	$\text{sum}(\text{dbinom}(a+1:b, \text{size}=n, \text{prob}=p))$
$P(x \geq a) \text{ o } 1 - P(x=0)$	$\text{sum}(\text{dbinom}(a:n, \text{size}=n, \text{prop}=p))$
$P(x > a)$	$\text{sum}(\text{dbinom}(a+1:n, \text{size}=n, \text{prop}=p))$

$P(z > a) \text{ o } 1 - P(z \leq a)$	$1 - \text{pnorm}(a, \text{mean}=\mu, \text{sd}=\sigma)$ $\text{pnorm}(a, \text{mean}=\mu, \text{sd}=\sigma, \text{lower.tail}=\text{FALSE})$
$P(a < Z < b)$	$\text{pnorm}(b, \text{mean}=\mu, \text{sd}=\sigma) - \text{pnorm}(a, \text{mean}=\mu, \text{sd}=\sigma)$
$P(x < a)$	$\text{pnorm}(a, \text{mean}=\mu, \text{sd}=\sigma)$

DISTRIB.	DENSITA' pdf	RIPARTIZIONE cdf	QUANTILI ORDINE α
$x \sim \text{unif}[a, b]$ uniforme	$\text{dunif}(x, \text{min}=a, \text{max}=b)$	$\text{punif}(x, \text{min}=a, \text{max}=b)$	$\text{qunif}(\alpha, \text{min}=a, \text{max}=b)$
$x \sim N(\mu, \sigma^2)$ normale standard	$\text{dnorm}(x, \text{mean}=\mu, \text{sd}=\sigma)$	$\text{pnorm}(x, \text{mean}=\mu, \text{sd}=\sigma)$	$\text{qnorm}(\alpha, \text{mean}=\mu, \text{sd}=\sigma)$
$x \sim B(n, p)$ quantili	$\text{dbinom}(x, \text{size}=n, \text{prob}=p)$	$\text{pbinom}(x, \text{size}=n, \text{prob}=p)$	$\text{qbinom}(\alpha, \text{size}=n, \text{prob}=p)$

❖ Statistica inferenziale

Media μ e deviazione standard σ NOTE

Valore atteso	$E[x] = \mu$
Varianza	$Var(x) = \sigma^2/n$
Teorema del limite centrale	$x \sim N(\mu, \sigma^2/n)$ con $n \geq 30$ quindi non necessaria normalità del campione
Errore statistico ($E > 0$)	$E = Z^* \sigma / \sqrt{n}$ con $Z^* = Z_{1-\frac{\alpha}{2}}$ <code>qnorm(1-(alpha/2), mean=0, sd=1)</code>
Livello di fiducia	$1-\alpha = \% \rightarrow \alpha = 1-cl$
Intervallo di confidenza, cl	$\mu = \bar{x} \pm E$ $IC = (\bar{x} - E, \bar{x} + E)$ $IC < -\bar{x} + c(-E, E)$

Media μ e varianza σ^2 NON NOTE, distr norm

Valore atteso	$E[x] = 0$
Varianza	$Var(x) = n/n-2$
t di student	t_{n-1}
Errore statistico	$E = t^* s / \sqrt{n}$ con $t^* = t_{1-\frac{\alpha}{2}, n-1}$ <code>qt(1-alpha/2, df=n-1)</code> se $n > 30$ $t_{1-\frac{\alpha}{2}, n-1}$ sostituito con $Z_{1-\frac{\alpha}{2}}$
Livello di fiducia	$1-\alpha/2$ con n-1 gradi di libertà
Intervallo di confidenza, cl	$\mu = \bar{x} \pm E$ $IC = (\bar{x} - E, \bar{x} + E)$ $IC < -\bar{x} + c(-E, E)$

Proporzione p di una popolazione bernoulliana

Proporzione campionaria, phat	$\text{phat} = n^\circ \text{successi nel campione} / n$
Ipotesi di lavoro	$n \cdot \text{phat} \geq 5$ e $n \cdot (1-\text{phat}) \geq 5$
Errore statistico	$E = Z_{1-\frac{\alpha}{2}} * \sqrt{\frac{\text{phat} \cdot (1-\text{phat})}{n}}$ <code>qnorm(1-(alpha/2), mean=0, sd=1) ↔ qnorm(1-(alpha/2))</code>
Livello di fiducia	$1-\alpha = \% \rightarrow \alpha = 1-cl$
Intervallo di confidenza, cl	$\mu = \text{phat} \pm E$ $IC = (\text{phat} - E, \text{phat} + E)$ $IC < -\text{phat} + c(-E, E)$

Stima della varianza, distr norm e varianza nota

Distribuzione chi-quadro con n gradi di libertà	<code>dchisq(x, df=n)</code>
Quantili	$l^* = X^2_{\alpha/2, n}$ <code>qchisq(alpha/2, df=n)</code> $r^* = X^2_{1-\alpha/2, n}$ <code>qchisq(1-alpha/2, df=n)</code>

Stima della varianza, distr norm e varianza NON nota

Livello di fiducia	$1-\alpha$
Quantili	$l^* = X^2_{\alpha/2, n-1}$ <code>qchisq(alpha/2, df=n-1)</code> $r^* = X^2_{1-\alpha/2, n-1}$ <code>qchisq(1-alpha/2, df=n-1)</code> <code>qchisq(alpha/2, df=n-1, lower.tail = FALSE)</code>
Intervallo di confidenza	$\left(\frac{(n-1)s^2}{r^*}, \frac{(n-1)s^2}{l^*} \right)$ $IC < -(n-1) * s^2 * c(1/rstar, 1/lstar)$

Livelli di confidenza minori forniscono stime intervallari più precise

TEST di ipotesi con popolazione bernoulliana $Z = \frac{p-p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$

Test a una cosa destra prop.test(x,n,p=p0,alternative="greater") con x=n° successi del campione, round(n*phat) n=ampiezza del campione p0=parametro teorico di confronto	H0:p=p0 HA:p>p0 Pvalue=P(Z>z)
Test a una coda sinistra prop.test(x,n,p=p0,alternative="less")	H0:p=p0 HA:p<p0 Pvalue=P(Z<z)
Test a 2 code prop.test(x,n,p=p0,alternative="two.sided")	H0:p=p0 HA:p≠p0 Pvalue=2P(Z> z)=2(1-P(Z< z))

TEST di ipotesi per media con varianza NOTA $Z = \frac{\bar{x}-\mu_0}{\frac{\sigma}{\sqrt{n}}}$

Test a una cosa destra ➤ pvalue<-1-pnorm(z) ➤ pvalue<-1-pnorm(z,mean=0,sd=1,lower.tail=FALSE)	H0:μ=μ0 HA: μ>μ0 Pvalue=P(Z>z)
Test a una coda sinistra pvalue<-pnorm(z,mean=0,sd=1)	H0:μ=μ0 HA: μ<μ0 Pvalue=P(Z<z)
Test a 2 code ➤ pvalue<-2*(1-pnorm(abs(z))) ➤ pvalue<-2*pt(abs(z),lower.tail=FALSE)	H0:μ=μ0 HA: μ ≠ μ0 Pvalue=2P(Z> z)=2(1-P(Z< z))

TEST di ipotesi per media con varianza NON NOTA $t = \frac{\bar{x}-\mu_0}{\frac{s}{\sqrt{n}}}$

Test a una cosa destra t.test(x,mu=mu0,alternative="greater")	H0:μ=μ0 HA: μ>μ0 Pvalue=P(T>t)
Test a una coda sinistra ➤ t.test(x,mu=mu0,alternative="less") ➤ pvalue<-pt(t,df=(n-1))	H0:μ=μ0 HA: μ<μ0 Pvalue=P(T<t)
Test a 2 code ➤ t.test(x,mu=mu0,alternative="two.sided") ➤ pvalue<-2*pt(abs(t),df=(n-1),lower.tail=FALSE) ➤ pvalue<-2*(1-pt(abs(t),df=(n-1)))	H0:μ=μ0 HA: μ ≠ μ0 Pvalue=2P(T> t)=2(1-P(T< t))

TEST di ipotesi per la varianza $\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$

Test a una cosa destra ➤ pvalue<-1-pchisq(chi2,df=n-1) ➤ pvalue<-pchisq(chi2,df=n-1, lower.tail=FALSE)	H0:σ ² = σ0 ² HA:σ ² > σ0 ² Pvalue=P(X>x ²)
Test a una coda sinistra pvalue<-pchisq(chi2,df=n-1)	H0:σ ² = σ0 ² HA:σ ² < σ0 ² Pvalue=P(X<x ²)
Test a 2 code pvalue<-2*min(pchisq(chi,df=(n-1),lower.tail=TRUE), pchisq(chi, df=(n-1),lower.tail=FALSE))	H0:σ ² = σ0 ² HA:σ ² ≠ σ0 ² Pvalue=2min(P(X<x ²);P(X>x ²))

Test di ipotesi per la mediana, distrib NON nota = TEST DI WILCOXON

Test a una cosa destra wilcox.test(x,mu=m,alternative="greater")	H0:mediana=m HA:mediana>m
Test a una coda sinistra wilcox.test(x,mu=m,alternative="less")	H0:mediana=m HA:mediana<m
Test a 2 code wilcox.test(x,mu=m,alternative="two.sided")	H0:mediana=m HA:mediana≠m

❖ Test di confronto tra 2 proporzioni di 2 popolazioni bernoulliane

$$z = \frac{p1 - p2}{\sqrt{\frac{p(1-p)}{n1 + n2}}}$$

z test

❖ Test di confronto tra medie di 2 popolazioni
-varianza note

$$z = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{\sigma_x^2 + \sigma_y^2}{n_x + n_y}}}$$

z test, con pvalue<-pnorm(z)

-varianze NON note, ma uguali

$$t = \frac{\bar{x} - \bar{y}}{sp \sqrt{\frac{1}{n_x + n_y}}} \quad sp = \text{deviazione standard pooled}$$

t.test(x,y,alternative="grater",var.equal=TRUE)
"less",
"two.sided",

-varianze NON note, NON uguali

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s_x^2 + s_y^2}{n_x + n_y}}} \quad f \text{ ha distribuzione di Fischer co } nA - 1 \text{ e } nB - 1 \text{ gradi di libert\`a}$$

t.test(x,y,alternative="grater",var.equal=FALSE)
"less",
"two.sided",

❖ Test di confronto sulle varianze

$$f = \frac{s^2_A}{s^2_B}$$

var.test(A,B,alternative="grater")
"less")
"two.sided")

❖ Test di confronto tra mediane

Wilcox.test(x,y,alternative="grater")
"less")
"two.sided")

TEST di indipendenza $n = (n^{\circ} \text{righe} - 1) * (n^{\circ} \text{colonne} - 1)$	H0: le variabili sono indipendenti HA: le variabili non sono indipendenti costruiamo la tabella di contingenza: <code>x<-rbind(rigasopra,rigasotto)</code> <code>chisq.test(x)</code>
TEST di Shapiro-Wilk	H0: X ha distribuzione normale HA: X NON ha distribuzione normale <code>shapiro.test(x)</code>
TEST di adattamento	H0: $X=X_0$ (modello compatibile con i dati) HA: $X \neq X_0$ (modello NON compatibile con i dati) <code>chisq.test(z,p=pt)</code> con <code>z=frequenze</code> , <code>pt=elementi</code>
TEST di Kolmogorov-Smirnov	H0: la distribuzione si adatta al modello HA: la distribuzione NON si adatta al modello <code>ks.test(x,"nome cdf",distribuzione)</code>
Confronto tra distribuzioni di 2 popolazioni	H0: $X=Y$ (distrib. 2 popolaz. sono uguali) HA: $X \neq Y$ (distrib. 2 popolaz. NON sono uguali) <code>ks.test(x,y)</code>
TEST ANOVA a una via distrib. normale e stesse varianze (anche non note)	H0: $\mu = \mu_0 = \dots = \mu_k$ HA: almeno una delle μ diversa dalle altre costruiamo un dataframe: <code>l<-list(maggio=may,settembre=sep,dicembre=dec)</code> <code>d<-stack(l)</code> <code>oneway.test(colonna1~colonna2,data=d,</code> <code>var.equal=TRUE)</code> <code>var.equal=FALSE)</code> → se non è ipotizzabile l'uguaglianza tra le varianze
TEST di Kruskal Wallis	H0: $\mu = \mu_0 = \dots = \mu_k$ HA: almeno una delle μ diversa dalle altre costruiamo un dataframe: <code>l<-list(maggio=may,settembre=sep,dicembre=dec)</code> <code>d<-stack(l)</code> <code>kruskal.test(colonna1~colonna2,data=d)</code>
TEST di ipotesi per (rho) $\rho=0$ Pvalue piccolo → t grande → r prossimo a 1 Pvalue grande → t piccolo → r prossimo a 0	H0: $\rho=0$ (NON c'è correlazione lineare) HA: $\rho \neq 0$ (c'è correlazione lineare) <code>Cor(x,y)</code> $t = r \left(\sqrt{\frac{n-2}{1-r^2}} \right) \quad n-2 \text{ gradi di libertà}$ <code>Pvalue=2P(T> t) → cor.test(x,y)</code>
TEST di pendenza e intercetta = 0	H0: $\beta_0=0$ intercetta ($\beta_1=0$ pendenza) HA: $\beta_0 \neq 0$ intercetta ($\beta_1 \neq 0$ pendenza) <code>reg<- lm(y~x)</code> <code>summary(reg) → summary(lm(y~x))</code>